

Deep Learning From Crowdsourced Labels: Coupled Cross-Entropy Minimization, Identifiability, and Regularization

Shahana Ibrahim, Tri Nguyen, Xiao Fu

ICLR 2023, Virtual Talk



Oregon State
University

Labeled Data

One of the pillars of AI revolution is *massive amount of labeled data*



Amazon team taps millions of Alexa interactions to reduce NLP error rate

KYLE WIGGERS @KYLE_L_WIGGERS
JANUARY 22, 2019 6:59 AM



Amazon VP of devices David Limp at a September 2018 event at Amazon headquarters in Seattle, Washington.

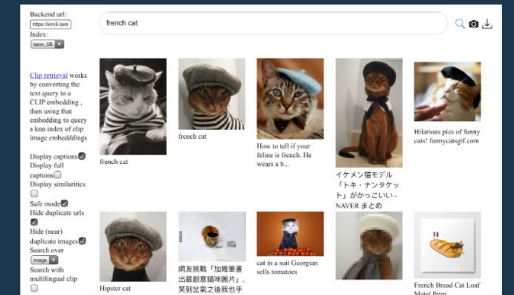
Image Credit: Khari Johnson / VentureBeat

LAION-5B: A NEW ERA OF OPEN LARGE-SCALE MULTI-MODAL DATASETS

by: Romain Beaumont, 31 Mar, 2022

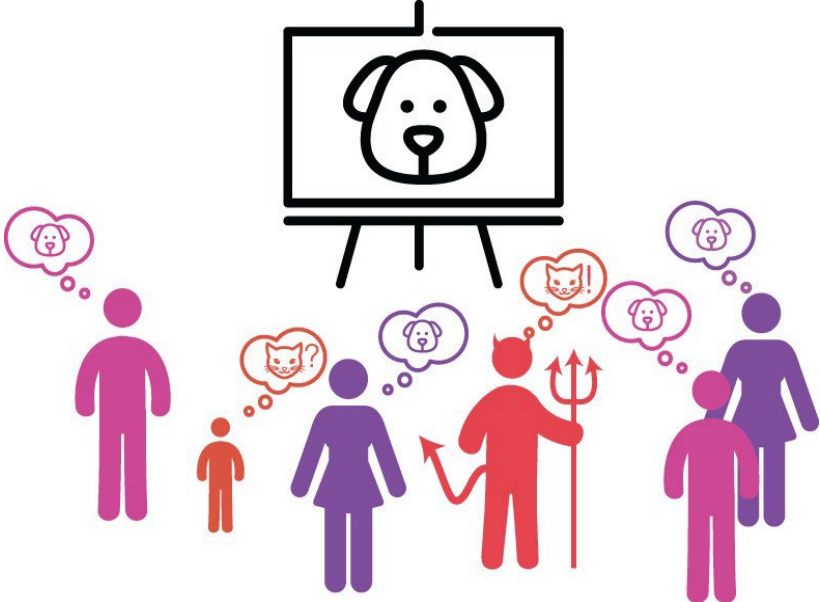
We present a dataset of 5.85 billion CLIP-filtered image-text pairs, 14x bigger than LAION-400M, previously the biggest openly accessible image-text dataset in the world - see also our [NeurIPS2022 paper](#)

Authors: Christoph Schuhmann, Richard Vencu, Romain Beaumont, Theo Coombes, Cade Gordon, Aarush Katta, Robert Kaczmarczyk, Jenia Jitsev



<https://laion.ai/blog/laion-5b/>

Crowdsourcing – Using Wisdom of the Crowd



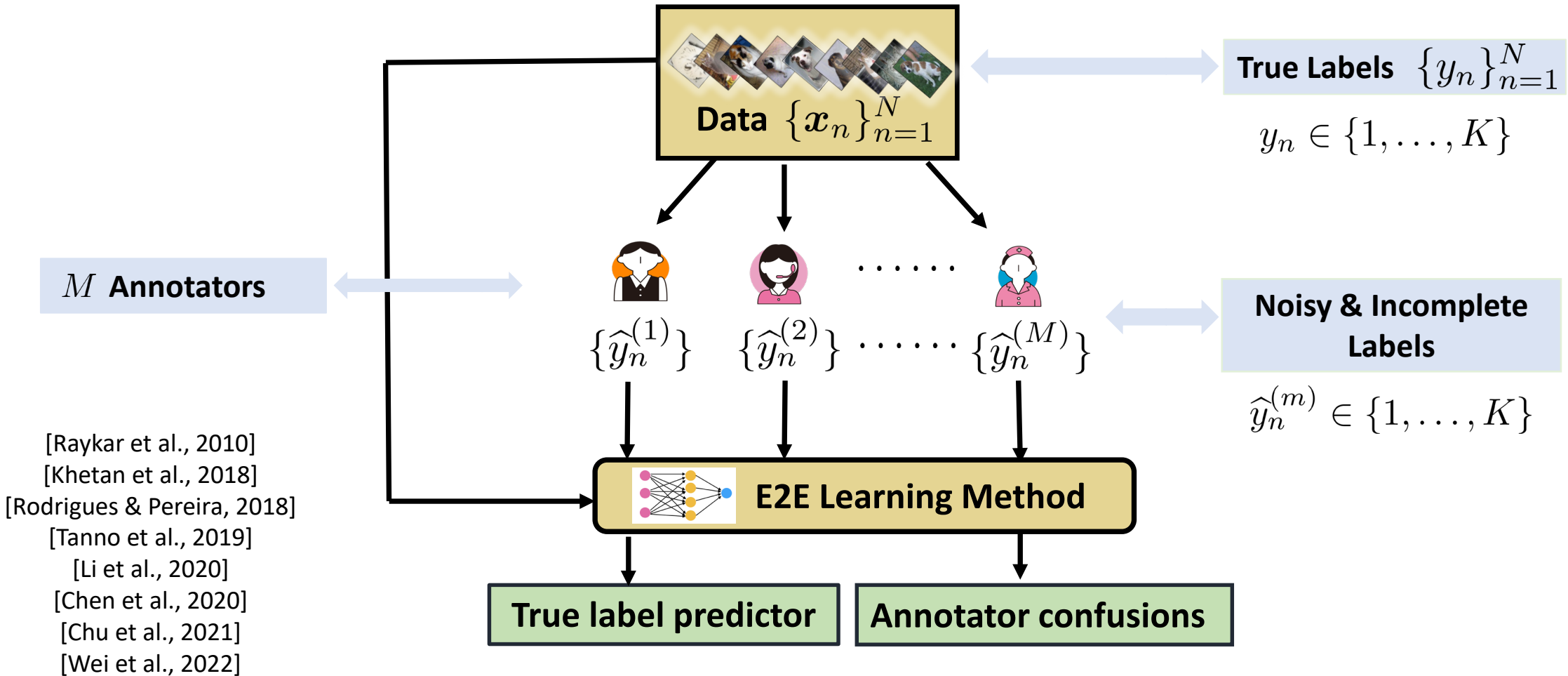
Each data item is labeled by multiple, often non-expert, annotators

→

Data	Labels			
	Cat	?	Dog
	Dog	Cat	?
⋮	⋮	⋮	⋮	⋮
	?	?	Dog

Noisy & Incomplete

End-To-End Learning for Crowdsourced Labels



Noisy Label Generation Model

$$\Pr(\hat{y}_n^{(m)} = k | \mathbf{x}_n) = \sum_{k'=1}^K \Pr(\hat{y}_n^{(m)} = k | y_n = k') \Pr(y_n = k' | \mathbf{x}_n)$$

Prob. Of m th
annotator's
response given \mathbf{x}_n

$\mathbf{p}_n^{(m)}$

Confusion
Matrix

\mathbf{A}_m

True label
predictor

$f(\mathbf{x}_n)$

$$\mathbf{p}_n^{(m)} = \mathbf{A}_m f(\mathbf{x}_n), \forall m, n$$

$$\hat{y}_n^{(m)} \sim \text{categorical}(\mathbf{p}_n^{(m)})$$

Goal : learn f and $\mathbf{A}_m, \forall m$

Assumption: Annotator
confusions are not data
dependent



Coupled Cross Entropy Minimization (CCEM)

The most popular E2E learning criterion [Rodrigues & Pereira, 2018]

$$\begin{aligned} & \text{minimize}_{f, \{A_m\}} \frac{1}{|\mathcal{S}|} \sum_{(m,n) \in \mathcal{S}} \text{CE}(A_m f(x_n), \hat{y}_n^{(m)}) \\ & \text{subject to } f \in \mathcal{F}, A_m \in \mathcal{A}, \forall m. \end{aligned}$$

[Rodrigues & Pereira, 2018]
[Tanno et al., 2019]
[Li et al., 2020]
[Chen et al., 2020]
[Chu et al., 2021]
[Wei et al., 2022]

Indices of
observed
labels

Neural
network
function class

Constrained set
 $\{A \in \mathbb{R}^{K \times K} \mid A \geq 0, \mathbf{1}^\top A = \mathbf{1}^\top\}$

Can CCEM learn the true label predictor and the true annotator confusions?

Analysis Result: CCEM correctly learns the true confusions and the true classifier,

under the assumptions

1 *Anchor point condition*



For each class, there is a data sample belonging to that class with prob. close to 1

2 *Class expert condition*



For each class, there is an expert which can predict that class correctly with prob. close to 1

Often, its hard to hold the two conditions together

Proposed Learning Criteria

✓ *Anchor point condition*

✗ *Class expert condition*



If we have more data,
but no experts to label

GeoCrowdNet (F)

$$\underset{\mathbf{f}, \{\mathbf{A}_m\}}{\text{minimize}} \quad \frac{1}{|\mathcal{S}|} \sum_{(m,n) \in \mathcal{S}} \text{CE}(\mathbf{A}_m \mathbf{f}(\mathbf{x}_n), \hat{y}_n^{(m)}) - \lambda \log \det \mathbf{F} \mathbf{F}^\top$$

subject to $\mathbf{f} \in \mathcal{F}, \mathbf{A}_m \in \mathcal{A}, \forall m.$

Regularization term

$$\begin{bmatrix} \mathbf{p}_1^{(1)} & \dots & \mathbf{p}_N^{(1)} \\ \vdots & \ddots & \vdots \\ \mathbf{p}_1^{(M)} & \dots & \mathbf{p}_N^{(M)} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_M \end{bmatrix} \begin{bmatrix} \mathbf{f}(\mathbf{x}_1) & \dots & \mathbf{f}(\mathbf{x}_N) \end{bmatrix}$$

\mathbf{P} \mathbf{W} \mathbf{F}

Underlying
NMF model

Proposed Learning Criteria

✗ *Anchor point condition*

✓ *Class expert condition*



If we have less data, but there are experts to label

GeoCrowdNet (W)

$$\underset{\mathbf{f}, \{\mathbf{A}_m\}}{\text{minimize}} \quad \frac{1}{|\mathcal{S}|} \sum_{(m,n) \in \mathcal{S}} \text{CE}(\mathbf{A}_m \mathbf{f}(\mathbf{x}_n), \hat{y}_n^{(m)}) - \lambda \log \det \mathbf{W}^\top \mathbf{W}$$

subject to $\mathbf{f} \in \mathcal{F}, \mathbf{A}_m \in \mathcal{A}, \forall m.$

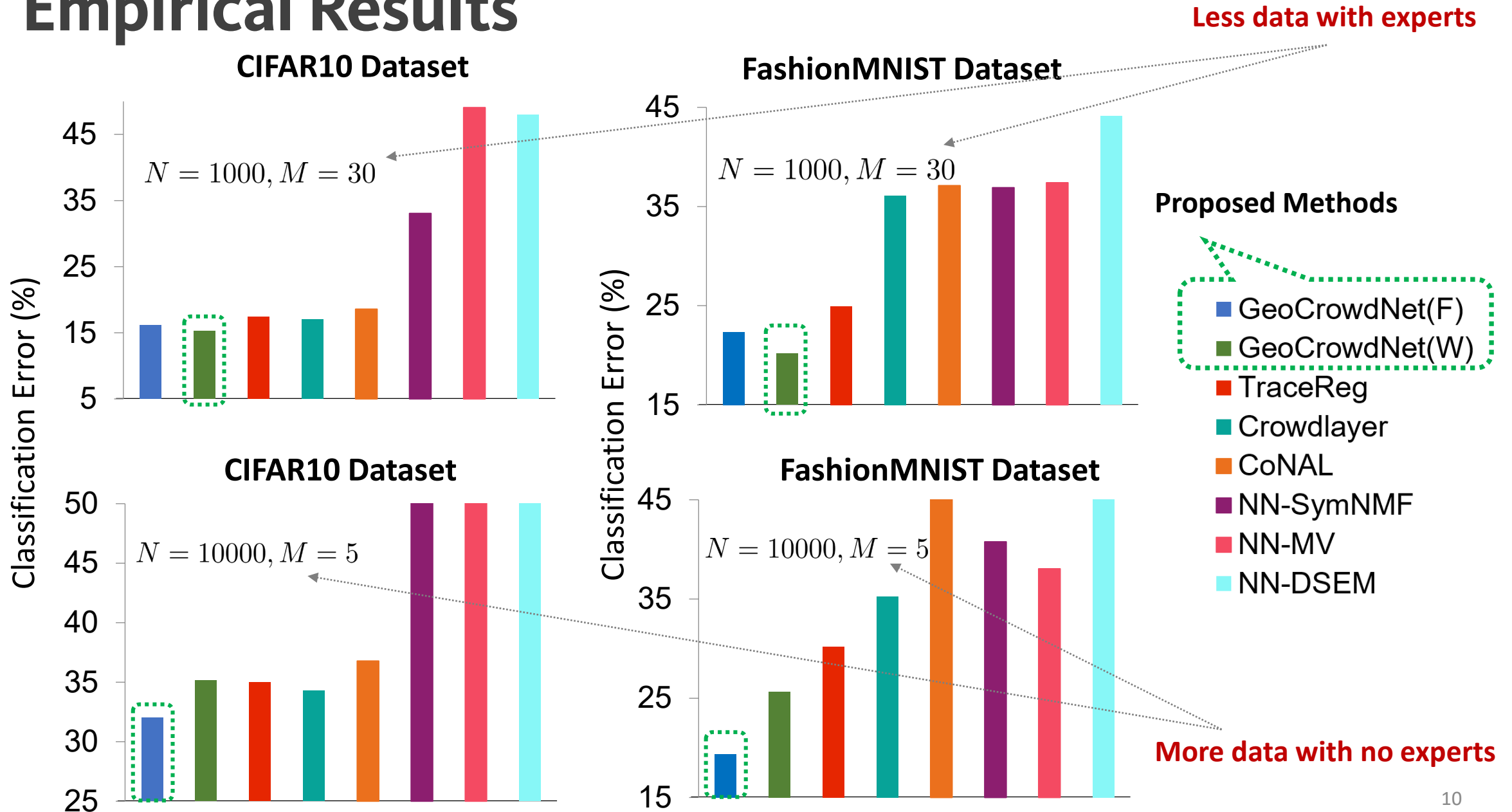
Regularization term

$$\begin{bmatrix} \mathbf{p}_1^{(1)} & \dots & \mathbf{p}_N^{(1)} \\ \vdots & \ddots & \vdots \\ \mathbf{p}_1^{(M)} & \dots & \mathbf{p}_N^{(M)} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_M \end{bmatrix} \begin{bmatrix} \mathbf{f}(\mathbf{x}_1) & \dots & \mathbf{f}(\mathbf{x}_N) \end{bmatrix}$$

\mathbf{P} \mathbf{W} \mathbf{F}

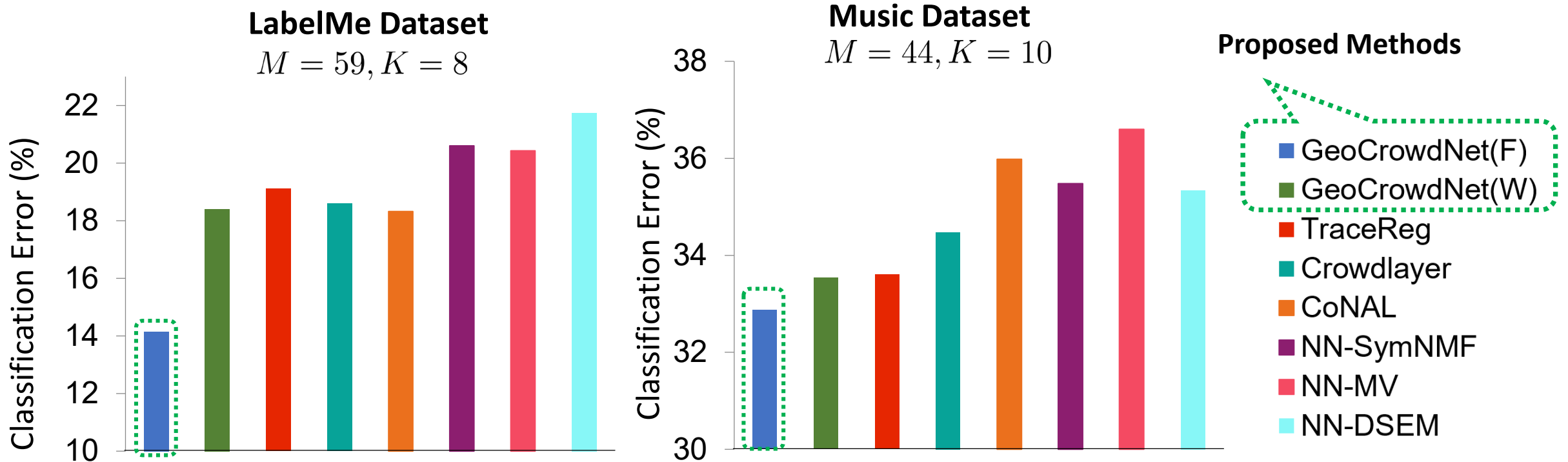
Underlying NMF model

Empirical Results



Empirical Results

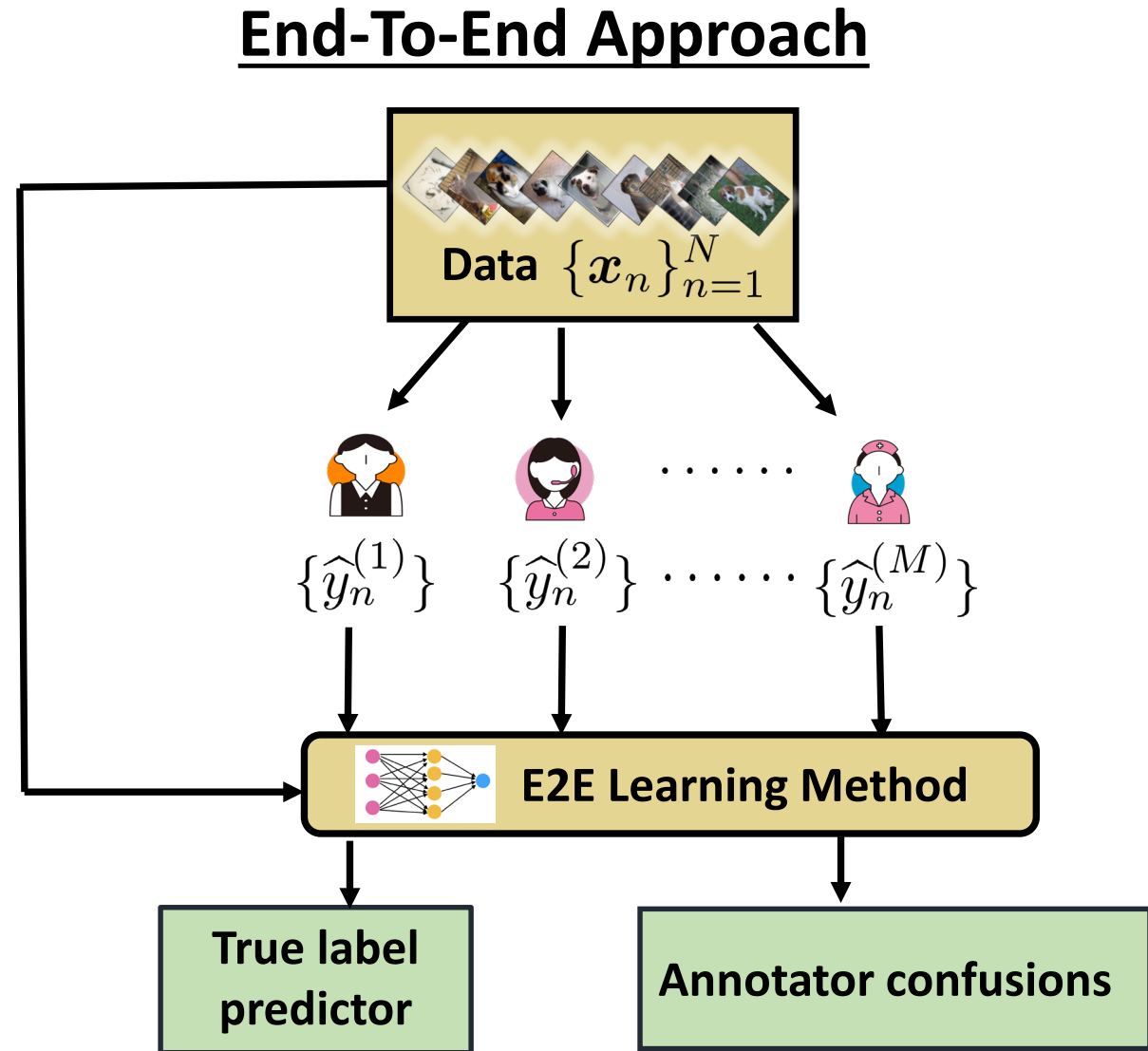
Noisy labels from  workers



The criterion designed for “no experts case” shows edge in practice

Key Takeaways

- ❑ Established deeper understanding to the challenging E2E learning problem
- ❑ Designed learning criteria with enhanced performance under practical scenarios, **e.g., no expert annotators**



Thank You